

CENTER FOR
BIOLOGICAL
CALSEQUENCE
ANALYSIS
LYSIS CBS

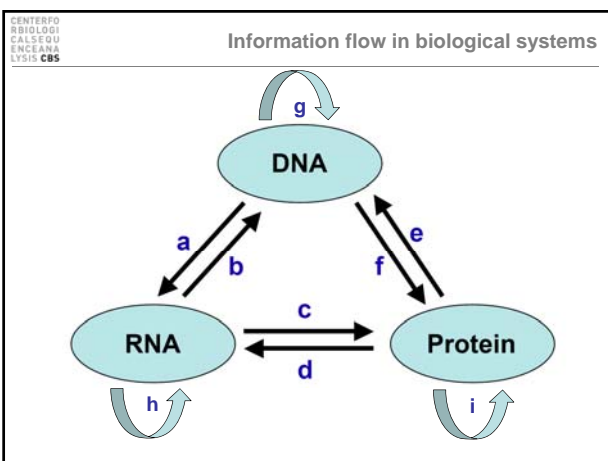
DNA as Biological Information

Rasmus Wernersson
Henrik Nielsen

CENTER FOR
BIOLOGICAL
CALSEQUENCE
ANALYSIS
LYSIS CBS

Overview

- Learning objectives
 - About Biological Information
 - A note about DNA sequencing techniques and DNA data
 - File formats used for biological data
 - Introduction to the GenBank database



DNA sequences = summary of information

5' AGCC 3'
3' TCGG 5'

5' ATGCCAGGTAA 3'

DNA backbone: <http://en.wikipedia.org/wiki/DNA>
(Deoxy)ribose: <http://en.wikipedia.org/>

**CENTER FOR
BIOLOGICAL
CALSEQU
ENGINEERING
ANALYSIS CBS**

PCR

Cycle 1

5'-CTGAGTATGAGACCTATAGGTACGGTGGCCATTCTGTCTGTGATCCCGGACTACTACAGAA-3'

|||||

3'-GGGCGATGATGG-5'

5'-ATGAGACCTATAG-3'

|||||

3'-GATCTTATCTTCTTGATATCCATGCCCGGATGAGTACGACTAGTGGCCATTGATGGATGCTT-5'

35 cycles

Melting
96°, 30 sec

↓

Annealing
~55°, 30 sec

↓

Extension
72°, 30 sec

CENTER FOR
BIOLOGICAL
CALSEQU
ENCEANA
LTSIO CBS

PCR

AMOUNT OF DNA

PCR CYCLE NUMBER

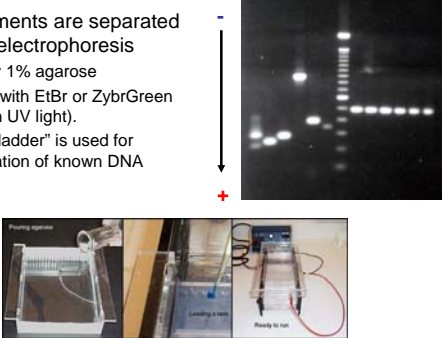
Real target

Single-primer target (500)

Single-primer target (1000)

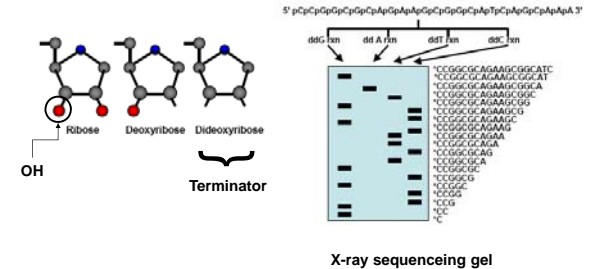
Gel electrophoresis

- DNA fragments are separated using gel electrophoresis
 - Typically 1% agarose
 - Colored with EtBr or ZybrGreen (glows in UV light).
 - A DNA "ladder" is used for identification of known DNA lengths.



Gel picture: <http://www.pharmaceutical-technology.com/projects/roche/images/roche3.jpg>
 PCR setup: <http://arbl.cvmbs.colostate.edu/bbooks/genetics/biotech/gels/agardna.html>

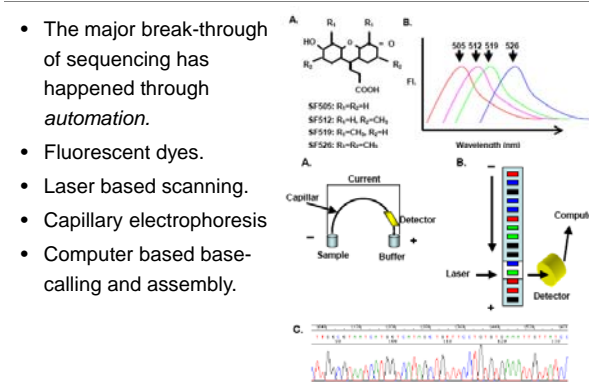
The Sanger method of DNA sequencing



Images: http://www.idtdna.com/support/technical/TechnicalBulletinPDF/DNA_Sequencing.pdf

Automated sequencing

- The major break-through of sequencing has happened through *automation*.
- Fluorescent dyes.
- Laser based scanning.
- Capillary electrophoresis
- Computer based base-calling and assembly.

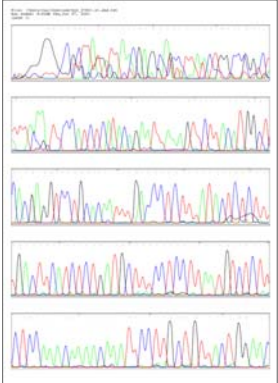


Images: http://www.idtdna.com/support/technical/TechnicalBulletinPDF/DNA_Sequencing.pdf

CENTER FOR
BIOLOGICAL
CALSEQU
ENGINEERING
LYSIS CBS

Handout exercise: "base-calling"

- Handout: Chromatogram
- Groups of 2-3.
- Tasks:
 - Identify "difficult" regions
 - Identify likely errors
 - Try to estimate the best interval to use



CENTER FOR
BIOLOGICAL
CALSEQU
ENGINEERING
LYSIS CBS

Sequence read mapping



CENTER FOR
BIOLOGICAL
CALSEQU
ENGINEERING
LYSIS CBS

DNA sequencing — history

1972 Recombinant DNA technology [Paul Berg].

1976 The first sequenced genome, the bacteriophage MS2 (actually, RNA) [Walter Fiers *et al.*]

1977 DNA sequencing by chemical cleavage [Allan Maxam & Walter Gilbert]; DNA sequencing by enzymatic synthesis [Fred Sanger].

1982 GenBank is established.

1987 The first automatic sequencer, *Prism 373* [Applied Biosystems].

1990 Human Genome Project is launched.

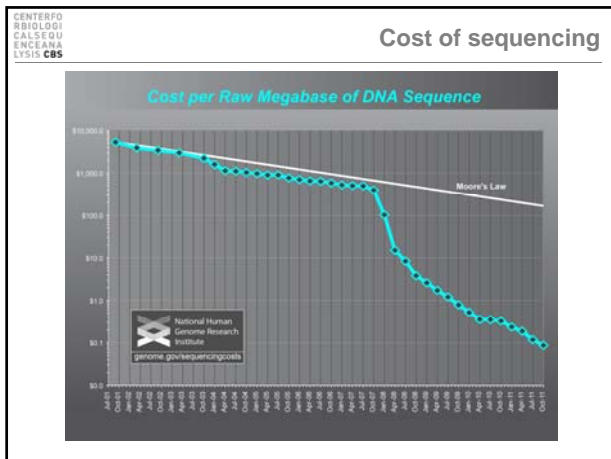
1995 The first genome of a free-living organism, the bacterium *Haemophilus influenzae* (1.8 Mb) [The Institute for Genomic Research (TIGR)].

1996 The first genome of a eukaryote, Baker's yeast, *Saccharomyces cerevisiae* (12.1 Mb) [International consortium].

1998 The first genome of an animal, the nematode *Caenorhabditis elegans* (97Mb) [Sanger Center *et al.*].

2001 The first "drafts" of the Human genome (3Gb) [Human Genome Project Consortium (Nature, 15 Feb) + Celera (Science, 16 Feb)].

Apr 15, 2012 GenBank release 189 contains 151,824,421 sequences with a total of 139,266,481,398 nucleotides (the files take up 586 GB).



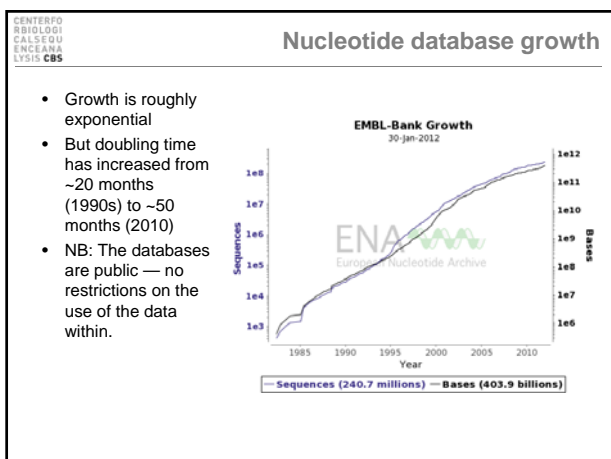
CENTER FOR BIOLOGICAL CALS EQU ENCEANA LYSIS CBS

Nucleotide databases

- **GenBank**, <http://www.ncbi.nlm.nih.gov/Genbank/>
- National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National Institutes of Health (NIH), USA
- Established in 1982.
- **EMBL**, <http://www.ebi.ac.uk/embl/>
- European Bioinformatics Institute (EBI), England
- Established in 1980 by the European Molecular Biology Laboratory, Heidelberg, Germany
- Now part of **ENA**, the European Nucleotide Archive, <http://www.ebi.ac.uk/ena/>
- **DDBJ**, <http://www.ddbj.nig.ac.jp/>
- National Institute of Genetics, Japan

Together they form

- International Nucleotide Sequence Database Collaboration, <http://www.insdc.org/>



CENTERFO
R BIOLOGI
CAL SEQU
ENCEANA
LYSIS CBS

FASTA format


```
>alpha-D
ATCTCGACCGACTCTGACAAAGAGCTGCTGCTGAGGTGTGGGGAAGGTGATCGGCCAC
CCAGACTGTGGAGCCGAGGCCCTGGAGAGGTGCGGGCTGAGCTTGGGGAAACCATGGGCA
AGGGGGGAGCTGGGTGGGAGCCCTACAGGGCTGCTGGGGGTGTTCGGCTGGGGGTGAG
CACTGACCATCCCGCTCCCGAGCTGTTCCACACTACCCCGAGACCAAGAGCTACTTCC
CCACTTCGACTTGCAACATGGCTCCGACAGGTCCGCAACCAAGCAAGAGGTGTGG
CCGCTTGGGCAAGCTGTCTAAGAGCTGGGCAACTCAAGCAAGCCCTGTCTGACTCA
GCAACTGCAAGCTCAACCTGCTGTGAGCCCTGTCAACTCAAGGCAAGCCGGGAGC
GGGGGTGAGGGGCGGGGAGTTGGGGGCGAGGAGCTGTTGGGGATCCGGGGCCATGCC
GGCGGTACTGAGCCCTGTTTGGCTTGCAGCTGTGGGCGAGTGTCTTCAAGTGGTGTG
GCCACACACTGGGCAAGACTACACCCGGAGGACATGCTGCTTTCGACAGTTCCTG
TGGGCTGTGTGACCGTGTGGCGAGAGTACAGATAA

>alpha-A
ATGCTGCTGTCTCCCAAGCAAGAGCACTGAAAGCCCTCTTCGCAAAATCGGCGGC
CAGGCCGAGTGAAGTGGGTGGAGCCCTGGAGAGTATGTGTCTATCCGTATTACCCC
ATCTCTTGTGTGTGTGTGATCCATCCATCTGCCCCATCTCTCCCATCCATACTG
TCCCTGTCTATGTGGCCCTGGCTCTGTCTCTGCTGCTGCTGCTGCTGCTGCTGCTG
TGTCCCGAGGTGTTCATCACTACCCCGAGCAAGAGCTACTTCCCGACTTGGACC
TGTCACTGGCTCCGCTCAGATCAAGGGGCAAGAGGTGGCGAGGCACTGGTTG
AGGCTGCCAAGCAAGTGAATGATCACTGCTGTGCTCTCCAGCTGAGGAGCTCCACG
CCCAAGGCTCGTGTGGAGCCCGCTCACTTCAAAGTGAAGTGGGAGGGGTGACCA
GTCTGGCTCCCGCTGACACACACTCTGGCTACCCCTGCACTCAACCCCTGTGTCACC
ATCTCTCTTGGCTTTCAGCTGTGGGTCACTGCTCTGCTGCTGCTGCTGCTGCTGCTG
CCCTCTCTCTGAGCCCGAGGTCCATGCTTCCCTGGCAAGTGTGTGTGCTGCTGGG
CACGCTCTTACTGCAAGTACGTTAA
```

(Handout)

CENTERFO
R BIOLOGI
CAL SEQU
ENCEANA
LYSIS CBS

GenBank format



- Originates from the GenBank database.
- Contains both a DNA sequence and annotation of feature (e.g. Location of genes).

(handout)

CENTERFO
R BIOLOGI
CAL SEQU
ENCEANA
LYSIS CBS

GenBank format - HEADER

```
LOCUS       CMGLOAD               1185 bp    DNA       linear   VRT 18-APR-2005
DEFINITION  Cairina moschata (duck) gene for alpha-D globin.
ACCESSION   X01831
VERSION     X01831.1  GI:62724
KEYWORDS    alpha-globin; globin.
SOURCE      Cairina moschata (Muscovy duck)
  ORGANISM  Cairina moschata
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Archosauria; Aves; Neognathae; Anseriformes; Anatidae; Cairina.
REFERENCE   1  (bases 1 to 1185)
AUTHORS     Erbil,C. and Niessing,J.
TITLE       The primary structure of the duck alpha D-globin gene: an unusual
            5' splice junction sequence
JOURNAL     EMBO J. 2 (8), 1339-1343 (1983)
PUBMED      10872328
COMMENT     Data kindly reviewed (13-NOV-1985) by J. Niessing.
```

Exercise: GenBank

GenBank format - FEATURE section

```
FEATURES             Location/Qualifiers
     source            1..1185
                        /organism="Cairina moschata"
                        /mol_type="genomic DNA"
                        /db_xref="taxon:8815"
CAAT_signal           20..24
                        /evidence=low
precursor_RNA         101..1114
                        /note="primary transcript"
exon                  101..234
                        /number=1
CDS                   join(143..234,387..591,939..1067)
                        /codon_start=1
                        /product="alpha D-glycin"
                        /protein_id="CAA25866.2"
                        /db_xref="GI:4453976"
                        /db_xref="GQA:P02003"
                        /db_xref="InterPro:IPR000971"
                        /db_xref="InterPro:IPR002338"
                        /db_xref="InterPro:IPR002340"
                        /db_xref="InterPro:IPR019050"
                        /db_xref="UniProt/Swiss-Prot:P02003"
                        /translation="M-L-T-A-R-W-I-V-L-S-G-E-N-G-A-D-F-P-S-R-L-G-M-G-L-Y-P-Q-T-K-T-P-P
F-I-F-L-E-P-S-Q-S-Q-W-O-S-K-V-K-A-A-A-G-N-A-G-W-E-L-N-S-G-A-S-E-L-S-H-L-S-A-T-N-L-R-V-P-F-P-K-I-L-L-A
Q-C-V-F-I-A-A-G-S-I-D-T-S-P-E-N-N-A-F-P-F-E-G-A-V-N-L-A-E-R"
repeat_region         227..246
                        /note="direct repeat 1"
intron                235..386
                        /number=1
repeat_region         288..308
                        /note="direct repeat 1"
exon                  387..591
                        /number=2
intron                592..939
                        /number=2
exon                  940..1114
                        /number=3
polyA_signal          1095..1100
                        1114
```

- Work in groups of 2-3 people.
- The exercise guide is linked from the course programme.
- Read the guide carefully - it contains a lot of information about GenBank.

